

Permutation Complexity and Coupling Measures in Hidden Markov Models

Taichi Haruna, Kohei Nakajima

Abstract—In [Haruna, T. and Nakajima, K., 2011. *Physica D* 240, 1370-1377], the authors introduced the duality between values (words) and orderings (permutations) as a basis to discuss the relationship between information theoretic measures for finite-alphabet stationary stochastic processes and their permutation versions. It has been used to give a simple proof of the equality between the entropy rate and the permutation entropy rate for any finite-alphabet stationary stochastic process and show some results on the excess entropy and the transfer entropy for finite-alphabet stationary ergodic Markov processes. In this paper, we generalize our previous framework and show the equalities between various information theoretic complexity and coupling measures and their permutation versions. In particular, we prove the following two results within the realm of hidden Markov models with ergodic internal processes: the two permutation versions of the transfer entropy, the symbolic transfer entropy and the transfer entropy on rank vectors, are both equivalent to the transfer entropy if they are considered as the rates, and the directed information theory can be captured by the permutation entropy approach.

Index Terms—Duality, Permutation Entropy, Excess Entropy, Transfer Entropy, Directed Information

I. INTRODUCTION

RECENTLY, the permutation-information theoretic approach to time series analysis proposed by Bandt and Pompe [1] has become popular in various fields [2]. It has been proved that the method of permutation is easy to implement relative to the other traditional methods, is computationally fast and is robust under the existence of noise [3], [4], [5]. However, if we turn our eyes to its theoretical side, few results are known for the permutation versions of information theoretic measures except the entropy rate.

There are two approaches to introduce permutation into dynamical systems theory. The first approach was introduced by Bandt et al. [6]. Given a one-dimensional interval map, they considered permutations induced by iterations of the map. Each point in the interval is classified into one of $n!$ permutations according to the permutation defined by $n - 1$ times iterations of the map starting from the point. Then, the Shannon entropy of this partition (called standard partition) of the interval is taken and normalized by n . The quantity obtained in the limit $n \rightarrow \infty$ is called permutation entropy if it exists. It was proved that the permutation entropy is equal to

the Kolmogorov-Sinai entropy for any piecewise monotone interval map [6]. This approach based on the standard partitions was extended by [7].

The second approach is taken by Amigó et al. [2], [8]. In this approach, given a measure-preserving map on a probability space, first an arbitrary finite partition of the space is taken. This gives rise to a finite-alphabet stationary stochastic process. An arbitrary ordering is introduced on the alphabet and the permutations of the words of finite lengths can be naturally defined (see Section II below). It is proved that the Shannon entropy of the occurrence of the permutations of a fixed length normalized by the length converges in the limit of the large length of the permutations. The quantity obtained is called permutation entropy rate (also called metric permutation entropy) and is shown to be equal to the entropy rate of the process. By taking the limit of finer partitions of the measurable space, the permutation entropy rate of the measure-preserving map is defined if the limit exists. Amigó [9] proved that it exists and is equal to the Kolmogorov-Sinai entropy.

In this paper, we restrict our attention to finite-alphabet stationary stochastic processes. Thus, we follow the second approach, namely, ordering on the alphabet is introduced arbitrarily. For quantities other than the entropy rate, three results for finite-alphabet stationary stochastic Markov processes have been shown by our previous work: the equality between the excess entropy and the permutation excess entropy [10], the equality between the mutual information expression of the excess entropy and its permutation version [11] and the equality between the transfer entropy rate and the symbolic transfer entropy rate [12].

The purpose of this paper is to set up a theoretical framework to discuss permutation versions of many information theoretic measures other than the entropy rate. In particular, we generalize our previous results for finite-alphabet stationary ergodic Markov processes to output processes of finite-state finite-alphabet hidden Markov models with ergodic internal processes. Upon this generalization, somewhat *ad hoc* proofs in our previous work become systematic and greatly simplified. This makes us easily access quantities that have not been considered in the permutation approach. In this paper, we shall treat the following quantities: excess entropy [13], transfer entropy [14], [15], momentary information transfer [16] and directed information [17], [18].

This paper is organized as follows: In Section II, we briefly review our previous result on the duality between words and permutations which is the basis for the succeeding results. In Section III, we prove a lemma about finite-state finite-alphabet hidden Markov models. In Section IV, we show

T. Haruna is with the Department of Earth & Planetary Sciences, Graduate School of Science, Kobe University, 1-1 Rokkodaicho, Nada, Kobe, 657-8501 Japan. e-mail: tharuna@penguin.kobe-u.ac.jp

K. Nakajima is with the Artificial Intelligence Laboratory, Department of Informatics, University of Zurich, Andreasstrasse 15, 8050 Zurich, Switzerland. e-mail: nakajima@ifi.uzh.ch

equalities between various information theoretic complexity and coupling measures and their permutation versions that hold for output processes of finite-state finite-alphabet hidden Markov models with ergodic internal processes. In Section V, we discuss how our results are related to the previous work in the literature.

II. THE DUALITY BETWEEN WORDS AND PERMUTATIONS

In this section, we summarize the results from our previous work [10] which will be used in this paper.

Let A_n be a finite set consisting of natural numbers from 1 to n called an *alphabet*. In this paper A_n is considered as a totally ordered set ordered by the usual ‘less-than-or-equal-to’ relationship. When we emphasize the total order, we call A_n *ordered alphabet*.

The set of all permutations of length $L \geq 1$ is denoted by \mathcal{S}_L . Namely, \mathcal{S}_L is the set of all bijections π on the set $\{1, 2, \dots, L\}$. For convenience, we sometimes denote a permutation π of length L by a string $\pi(1)\pi(2)\dots\pi(L)$. The number of *descents*, places with $\pi(i) > \pi(i+1)$, of $\pi \in \mathcal{S}_L$ is denoted by $\text{Desc}(\pi)$. For example, if $\pi \in \mathcal{S}_5$ is given by $\pi(1)\pi(2)\pi(3)\pi(4)\pi(5) = 35142$, then $\text{Desc}(\pi) = 2$.

Let $A_n^L = \underbrace{A_n \times \dots \times A_n}_L$ be the L -fold product of A_n .

A word of length $L \geq 1$ is an element of A_n^L . It is denoted by $x_{1:L} := x_1 \dots x_L := (x_1, \dots, x_L) \in A_n^L$. We say that the *permutation type* of a word $x_{1:L}$ is $\pi \in \mathcal{S}_L$ if we have $x_{\pi(i)} \leq x_{\pi(i+1)}$ and $\pi(i) < \pi(i+1)$ when $x_{\pi(i)} = x_{\pi(i+1)}$ for $i = 1, 2, \dots, L-1$. Namely, the permutation type of $x_{1:L}$ is the permutation of indices defined by re-ordering symbols x_1, \dots, x_L in the increasing order. For example, the permutation type of $x_{1:5} = 31212 \in A_3^5$ is $\pi(1)\pi(2)\pi(3)\pi(4)\pi(5) = 24351$ because $x_2x_4x_3x_5x_1 = 11223$.

Let $\phi_{n,L} : A_n^L \rightarrow \mathcal{S}_L$ be a map sending each word $x_{1:L}$ to its permutation type $\pi = \phi_{n,L}(x_{1:L})$. We define another map $\mu_{n,L} : \phi_{n,L}(A_n^L) \subseteq \mathcal{S}_L \rightarrow A_n^L$ by the following procedure:

- (i) Given a permutation $\pi \in \phi_{n,L}(A_n^L) \subseteq \mathcal{S}_L$, we decompose the sequence $\pi(1) \dots \pi(L)$ of length L into *maximal ascending subsequences*. A subsequence $i_j \dots i_{j+k}$ of a sequence $i_1 \dots i_L$ of length L is called a *maximal ascending subsequence* if it is ascending, namely, $i_j \leq i_{j+1} \leq \dots \leq i_{j+k}$, and neither $i_{j-1}i_j \dots i_{j+k}$ nor $i_ji_{j+1} \dots i_{j+k+1}$ is ascending.
- (ii) If $\pi(1) \dots \pi(i_1), \pi(i_1+1) \dots \pi(i_2), \dots, \pi(i_{k-1}+1) \dots \pi(L)$ is a decomposition of $\pi(1) \dots \pi(L)$ into maximal ascending subsequences, then a word $x_{1:L} \in A_n^L$ is defined by $x_{\pi(1)} = \dots = x_{\pi(i_1)} = 1, x_{\pi(i_1+1)} = \dots = x_{\pi(i_2)} = 2, \dots, x_{\pi(i_{k-1}+1)} = \dots = x_{\pi(L)} = k$. We define $\mu_{n,L}(\pi) = x_{1:L}$. Note that $\text{Desc}(\pi) \leq n-1$ because π is the permutation type of some word $y_{1:L} \in A_n^L$. Thus, we have $k = \text{Desc}(\pi) + 1 \leq n$. Hence, $\mu_{n,L}$ is well-defined as a map from $\phi_{n,L}(A_n^L)$ to A_n^L .

By construction, we have $\phi_{n,L} \circ \mu_{n,L}(\pi) = \pi$ for all $\pi \in \phi_{n,L}(A_n^L)$. To illustrate the construction of $\mu_{n,L}$, let us consider a word $y_{1:5} = 21123 \in A_3^5$. The permutation type of $y_{1:5}$ is $\pi(1)\pi(2)\pi(3)\pi(4)\pi(5) = 23145$. The decomposition of 23145 into maximal ascending subsequences is 23, 145.

We obtain $\mu_{n,L}(\pi) = x_{1:5} = 21122$ by putting $x_2x_3x_1x_4x_5 = 11222$.

Theorem 1: (i) For any $\pi \in \mathcal{S}_L$,

$$|\phi_{n,L}^{-1}(\pi)| = \binom{L+n-\text{Desc}(\pi)-1}{L},$$

where $\binom{a}{b} = 0$ if $a < b$.

- (ii) Let us put $B_{n,L} := \{x_{1:L} \in A_n^L | \phi_{n,L}^{-1}(\pi) = \{x_{1:L}\} \text{ for some } \pi \in \mathcal{S}_L\}$ and $C_{n,L} := \{\pi \in \mathcal{S}_L | |\phi_{n,L}^{-1}(\pi)| = 1\}$. Then, $\phi_{n,L}$ restricted on $B_{n,L}$ is a map into $C_{n,L}$, $\mu_{n,L}$ restricted on $C_{n,L}$ is a map into $B_{n,L}$, and they form a pair of mutually inverse maps. Furthermore, we have $B_{n,L} = \{x_{1:L} \in A_n^L | 1 \leq \forall i \leq n-1 \leq \exists j < k \leq L \text{ s. t. } x_j = i+1, x_k = i\}$ and $C_{n,L} = \{\pi \in \mathcal{S}_L | \text{Desc}(\pi) = n-1\}$.

Proof: The theorem is a recasting of statements in Lemma 5 and Theorem 9 in [10]. ■

Let $\mathbf{X} = \{X_1, X_2, \dots\}$ be a finite-alphabet stationary stochastic process, where each stochastic variable X_i takes its value in A_n . By the assumed stationarity, the probability of the occurrence of any word $x_{1:L} \in A_n^L$ is time-shift invariant: $\Pr\{X_1 = x_1, \dots, X_L = x_L\} = \Pr\{X_{k+1} = x_1, \dots, X_{k+L} = x_L\}$ for any $k, L \geq 1$. Hence, it makes sense to define it without referring to the time to start. We denote the probability of the occurrence of a word $x_{1:L} \in A_n^L$ by $p(x_{1:L}) = p(x_1 \dots x_L)$. The probability of the occurrence of a permutation $\pi \in \mathcal{S}_L$ is given by $p(\pi) = \sum_{x_{1:L} \in \phi_{n,L}^{-1}(\pi)} p(x_{1:L})$.

For a finite-alphabet stationary stochastic process \mathbf{X} over the alphabet A_n , we define

$$\alpha_{\mathbf{X},L} := \sum_{\substack{\pi \in \mathcal{S}_L, \\ |\phi_{n,L}^{-1}(\pi)| > 1}} p(\pi) = \sum_{\pi \notin C_{n,L}} p(\pi)$$

and

$$\begin{aligned} \beta_{x,\mathbf{X},L} &= \Pr\{x_{1:N} \in A_n^N | x_j \neq x \text{ for any } 1 \leq j \leq N\} \\ &= \sum_{\substack{x_j \neq x, \\ 1 \leq j \leq N}} p(x_1 \dots x_N), \end{aligned}$$

where $L \geq 1, x \in A_n$ and $N = \lfloor L/2 \rfloor$ and $\lfloor a \rfloor$ is the largest integer not greater than a .

Lemma 2: Let \mathbf{X} be a finite-alphabet stationary stochastic process and ϵ be a positive real number. If $\beta_{x,\mathbf{X},L} < \epsilon$ for any $x \in A_n$, then we have $\alpha_{\mathbf{X},L} < 2n\epsilon$.

Proof: The claim follows from Theorem 1 (ii). See Lemma 12 in [10] for the complete proof. ■

III. A RESULT ON FINITE-STATE FINITE-ALPHABET HIDDEN MARKOV MODELS

A *finite-state finite-alphabet hidden Markov model* (in short, HMM) [19] is a quadruple $(\Sigma, A, \{T^{(a)}\}_{a \in A}, \mu)$, where Σ and A are finite sets called *state set* and *alphabet*, respectively, $\{T^{(a)}\}_{a \in A}$ is a family of $|\Sigma| \times |\Sigma|$ matrices indexed by elements of A where $|\Sigma|$ is the size of state set Σ , and μ is a probability distribution on the set Σ . The following conditions must be satisfied:

- (i) $T_{ss'}^{(a)} \geq 0$ for any $s, s' \in \Sigma$ and $a \in A$,
- (ii) $\sum_{s', a} T_{ss'}^{(a)} = 1$ for any $s \in \Sigma$,
- (iii) and $\mu(s') = \sum_{s, a} \mu(s) T_{ss'}^{(a)}$ for any $s' \in \Sigma$.

Any probability distribution satisfying the condition (iii) is called a *stationary distribution*. The $|\Sigma| \times |\Sigma|$ matrix $T := \sum_{a \in A} T^{(a)}$ is called *state transition matrix*. The ternary (Σ, T, μ) defines the *underlying Markov chain*. Note that the condition (iii) is equivalent to the condition (iii') $\mu(s') = \sum_s \mu(s) T_{ss'}$.

Two finite-alphabet stationary processes are induced by a HMM $(\Sigma, A, \{T^{(a)}\}_{a \in A}, \mu)$. One is solely determined by the underlying Markov chain. It is called *internal process* and is denoted by $\mathbf{S} = \{S_1, S_2, \dots\}$. The alphabet for \mathbf{S} is Σ . The joint probability distributions which characterize \mathbf{S} is given by $\Pr\{S_1 = s_1, S_2 = s_2, \dots, S_L = s_L\} := \mu(s_1) T_{s_1 s_2} \dots T_{s_{L-1} s_L}$ for any $s_1, \dots, s_L \in \Sigma$ and $L \geq 1$. The other process $\mathbf{X} = \{X_1, X_2, \dots\}$ with the alphabet A is defined by the joint probability distributions $\Pr\{X_1 = x_1, X_2 = x_2, \dots, X_L = x_L\} := \sum_{s, s'} \mu(s) (T^{(x_1)} \dots T^{(x_L)})_{ss'}$ for any $x_1, \dots, x_L \in A$ and $L \geq 1$ and called *output process*. The stationarity of the probability distribution μ ensures that of both the internal and output processes.

Symbols $a \in A$ such that $T^{(a)} = O$ occur in the output process with probability 0. Hence, we obtain the same output process even if we remove these symbols. Thus, we can assume $T^{(a)} \neq O$ for any $a \in A$ without loss of generality.

The internal process \mathbf{S} of a HMM $(\Sigma, A, \{T^{(a)}\}_{a \in A}, \mu)$ is called *ergodic* if the state transition matrix T is *irreducible* [20]: for any $s, s' \in \Sigma$ there exists $k > 0$ such that $(T^k)_{ss'} > 0$. If the internal process \mathbf{S} is ergodic, then the stationary distribution μ is uniquely determined by the state transition matrix T via the condition (iii'). It is known that the ergodicity of the internal process \mathbf{S} implies that of the output process \mathbf{X} , but not vice versa [21].

Note that there are two types of hidden Markov models depending on whether outputs are emitted from edges or states. The HMM defined here is edge emitting type. However, it is known that these two classes of HMM are equivalent [19]. In particular, any finite-alphabet finite-order stationary Markov process can be described as a HMM defined here.

Lemma 3: Let \mathbf{X} be the output process of a HMM $(\Sigma, A_n, \{T^{(a)}\}_{a \in A_n}, \mu)$, where $A_n = \{1, 2, \dots, n\}$ is an ordered alphabet. If the internal process \mathbf{S} of the HMM is ergodic, then for any $x \in A_n$ there exists $0 < \gamma_x < 1$ and $C_x > 0$ such that $\beta_{x, \mathbf{X}, L} < C_x \gamma_x^L$ for any $L \geq 1$.

Proof: Given $L \geq 1$, let us put $N := \lfloor L/2 \rfloor$. Fix any $x \in A_n$. Since we have

$$\begin{aligned} \beta_{x, \mathbf{X}, L} &= \sum_{\substack{x_j \neq x, \\ 1 \leq j \leq N}} p(x_1 \dots x_N) \\ &= \sum_{\substack{x_j \neq x, \\ 1 \leq j \leq N}} \sum_{s, s'} \mu(s) (T^{(x_1)} \dots T^{(x_N)})_{ss'} \\ &= \langle \mu (T - T^{(x)})^N, \mathbf{1} \rangle, \end{aligned}$$

where $\mathbf{1} = (1, 1, \dots, 1)$ and $\langle \dots, \dots \rangle$ is the usual inner product of the $|\Sigma|$ -dimensional Euclidean space $\mathbb{R}^{|\Sigma|}$, it is sufficient to show that the largest eigenvalue of the matrix $T_{(x)} := T - T^{(x)}$ is less than 1. To prove this we shall appeal to the Perron-Frobenius theorem because $T_{(x)}$ is a non-negative matrix:

- (i) there exists a non-negative eigenvalue λ called the *Perron-Frobenius eigenvalue* such that any other eigenvalue of $T_{(x)}$ has absolute value not greater than λ ,
- (ii) $\lambda \leq \max_s \{\sum_{s'} (T_{(x)})_{ss'}\} \leq 1$,
- (iii) and there exists a non-negative left eigenvector \mathbf{v} corresponding to the eigenvalue λ .

We can show that for any $\epsilon > 0$ there exists $C_\epsilon > 0$ such that for any $k \geq 1$

$$\|\mu T_{(x)}^k\| \leq C_\epsilon (\lambda + \epsilon)^k \|\mu\|,$$

where $\|\dots\|$ is the Euclidean norm and we used the fact that any non-negative matrix and its transpose have the same Perron-Frobenius eigenvalue. For the proof of this inequality, see the beginning of section 1.2 in [22], for example. If $\lambda < 1$ then we can choose $\epsilon > 0$ so that $\lambda + \epsilon < 1$. If we put $\gamma_x := (\lambda + \epsilon)^{\frac{1}{2}}$ and $C_x := C_\epsilon (\lambda + \epsilon)^{-1} \|\mu\| \|\mathbf{1}\|$ then we obtain $\beta_{x, \mathbf{X}, L} < C_x \gamma_x^L$ by the Cauchy-Schwartz inequality as desired.

Let us derive a contradiction from the assumption $\lambda = 1$. If $\lambda = 1$ then we have $\mathbf{v} T_{(x)} = \mathbf{v}$. For any $k \geq 1$, We have

$$\langle \mathbf{v}, \mathbf{1} \rangle = \langle \mathbf{v} T_{(x)}^k, \mathbf{1} \rangle \leq \langle \mathbf{v} T^k, \mathbf{1} \rangle = \langle \mathbf{v}, T^k \mathbf{1} \rangle = \langle \mathbf{v}, \mathbf{1} \rangle,$$

because $T_{(x)} \leq T$ and T is a stochastic matrix. Thus, we obtain $\langle \mathbf{v} (T^k - T_{(x)}^k), \mathbf{1} \rangle = 0$. Since $\mathbf{1}$ is a positive vector and $\mathbf{v} (T^k - T_{(x)}^k)$ is a non-negative vector, it follows that

$$\mathbf{v} (T^k - T_{(x)}^k) = \mathbf{0}.$$

Let us consider $u, u' \in S$ such that $T_{uu'}^{(x)} > 0$. For any $s, s' \in S$, there exist $k_1, k_2 \geq 1$ such that $(T^{k_1})_{su} > 0$ and $(T^{k_2})_{u's'} > 0$ because T is irreducible. If we put $k = k_1 + k_2 + 1$ then it holds that

$$\begin{aligned} (T^k - T_{(x)}^k)_{ss'} &= \sum_{\substack{x_1, \dots, x_k, \\ \exists i \text{ s. t. } x_i = x}} (T^{(x_1)} \dots T^{(x_k)})_{ss'} \\ &\geq (T^{k_1})_{su} T_{uu'}^{(x)} (T^{k_2})_{u's'} > 0. \end{aligned}$$

On the other hand, the s' -th component of the vector $\mathbf{v} (T^k - T_{(x)}^k)$ must be 0:

$$\sum_{s''} v_{s''} (T^k - T_{(x)}^k)_{s'' s'} = 0,$$

where $v_{s''}$ denotes the s'' -th component of \mathbf{v} . We obtain $v_s = 0$ because \mathbf{v} is a non-negative vector and $(T^k - T_{(x)}^k)$ is a non-negative matrix. Since $s \in S$ is arbitrary, we conclude that $\mathbf{v} = \mathbf{0}$. However, this contradicts the fact that \mathbf{v} is an eigenvector. ■

IV. PERMUTATION COMPLEXITY AND COUPLING MEASURES

In this section, we discuss the equalities between complexity and coupling measures and their permutation versions for the output processes of HMMs whose internal processes are ergodic.

A. Fundamental lemma

Let $(\mathbf{X}^1, \dots, \mathbf{X}^m)$ be a multivariate finite-alphabet stationary stochastic process, where each univariate process $\mathbf{X}^k = \{X_1^k, X_2^k, \dots\}$, $k = 1, 2, \dots, m$ is defined over an ordered alphabet A_{n_k} . For simplicity, we use the notations

$$\begin{aligned} & p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m) \\ & := \Pr\{X_{a_1:b_1}^1 = x_{a_1:b_1}^1, \dots, X_{a_m:b_m}^m = x_{a_m:b_m}^m\}, \end{aligned}$$

$$\begin{aligned} & p(\pi_1, \dots, \pi_m) \\ & := \Pr\{\phi_{n_k, b_k - a_k + 1} \circ X_{a_k:b_k}^k = \pi_k, \quad k = 1, \dots, m\} \end{aligned}$$

and

$$p(\pi_k) := \Pr\{\phi_{n_k, b_k - a_k + 1} \circ X_{a_k:b_k}^k = \pi_k\},$$

where $1 \leq a_k \leq b_k$, $x_{a_k:b_k}^k \in A_{n_k}^{b_k - a_k + 1}$ and $\pi_k \in \mathcal{S}_{b_k - a_k + 1}$ for $k = 1, \dots, m$.

Lemma 4:

$$\begin{aligned} 0 & \leq H(X_{a_1:b_1}^1, \dots, X_{a_m:b_m}^m) - H^*(X_{a_1:b_1}^1, \dots, X_{a_m:b_m}^m) \\ & \leq \left(\sum_{k=1}^m \alpha_{\mathbf{X}^k, b_k - a_k + 1} \right) \left(\sum_{k=1}^m n_k \log(b_k - a_k + 1 + n_k) \right), \end{aligned}$$

where

$$\begin{aligned} & H(X_{a_1:b_1}^1, \dots, X_{a_m:b_m}^m) \\ & = - \sum_{x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m} p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m) \\ & \quad \times \log p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m) \end{aligned}$$

and

$$\begin{aligned} & H^*(X_{a_1:b_1}^1, \dots, X_{a_m:b_m}^m) \\ & = - \sum_{\pi_1, \dots, \pi_m} p(\pi_1, \dots, \pi_m) \log p(\pi_1, \dots, \pi_m) \end{aligned}$$

are the Shannon entropy of the joint occurrence of words $x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m$ and permutations π_1, \dots, π_m , respectively, and the base of the logarithm is taken as 2.

Proof: We have

$$\begin{aligned} & H(X_{a_1:b_1}^1, \dots, X_{a_m:b_m}^m) - H^*(X_{a_1:b_1}^1, \dots, X_{a_m:b_m}^m) \\ & = \sum_{\pi_1, \dots, \pi_m, p(\pi_1, \dots, \pi_m) > 0} p(\pi_1, \dots, \pi_m) \\ & \quad \times \left(- \sum_{\substack{x_{a_k:b_k}^k \in \phi_{n_k, b_k - a_k + 1}^{-1}(\pi_k), \\ 1 \leq k \leq m}} \frac{p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m)}{p(\pi_1, \dots, \pi_m)} \right. \\ & \quad \left. \times \log \frac{p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m)}{p(\pi_1, \dots, \pi_m)} \right). \end{aligned}$$

By Theorem 1 (i), it holds that

$$\begin{aligned} 0 & \leq - \sum_{\substack{x_{a_k:b_k}^k \in \phi_{n_k, b_k - a_k + 1}^{-1}(\pi_k), \\ 1 \leq k \leq m}} \frac{p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m)}{p(\pi_1, \dots, \pi_m)} \\ & \quad \times \log \frac{p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m)}{p(\pi_1, \dots, \pi_m)} \\ & \leq \log \prod_{k=1}^m \binom{b_k - a_k + n_k - \text{Desc}(\pi_k)}{b_k - a_k + 1} \\ & \leq \log \prod_{k=1}^m (b_k - a_k + 1 + n_k)^{n_k} \\ & = \sum_{k=1}^m n_k \log(b_k - a_k + 1 + n_k) \end{aligned}$$

for $(\pi_1, \dots, \pi_m) \in \mathcal{S}_{b_1 - a_1 + 1} \times \dots \times \mathcal{S}_{b_m - a_m + 1}$ such that $p(\pi_1, \dots, \pi_m) > 0$.

If $|(\phi_{n_1, b_1 - a_1 + 1} \times \dots \times \phi_{n_m, b_m - a_m + 1})^{-1}(\pi_1, \dots, \pi_m)| = 1$ then

$$\begin{aligned} & - \sum_{\substack{x_{a_k:b_k}^k \in \phi_{n_k, b_k - a_k + 1}^{-1}(\pi_k), \\ 1 \leq k \leq m}} \frac{p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m)}{p(\pi_1, \dots, \pi_m)} \\ & \quad \times \log \frac{p(x_{a_1:b_1}^1, \dots, x_{a_m:b_m}^m)}{p(\pi_1, \dots, \pi_m)} = 0. \end{aligned}$$

On the other hand, we have

$$\begin{aligned} & \sum_{\substack{\pi_1, \dots, \pi_m, \\ \exists k \text{ s.t. } |\phi_{n_k, b_k - a_k + 1}^{-1}(\pi_k)| > 1}} p(\pi_1, \dots, \pi_m) \\ & \leq \sum_{k=1}^m \sum_{\substack{\pi_k, \\ |\phi_{n_k, b_k - a_k + 1}^{-1}(\pi_k)| > 1}} p(\pi_k) \\ & = \sum_{k=1}^m \alpha_{\mathbf{X}^k, b_k - a_k + 1}. \end{aligned}$$

This completes the proof of the inequality. \blacksquare

B. Excess Entropy

Let \mathbf{X} be a finite-alphabet stationary stochastic process. Its *excess entropy* is defined by [13]

$$\begin{aligned} \mathbf{E}(\mathbf{X}) & = \lim_{L \rightarrow \infty} [H(X_{1:L}) - h(\mathbf{X})L] \\ & = \sum_{L=1}^{\infty} [H(X_L | X_{1:L-1}) - h(\mathbf{X})], \end{aligned}$$

if the limit on the right-hand side exists, where $h(\mathbf{X}) = \lim_{L \rightarrow \infty} H(X_{1:L})/L$ is the *entropy rate* of \mathbf{X} which exists for any finite-alphabet stationary stochastic process [23].

The excess entropy has been used as a measure of complexity [24], [25], [26], [27], [28], [29]. Actually, it quantifies global correlations present in a given stationary process in the following sense: if $\mathbf{E}(\mathbf{X})$ exists then it can be written as the mutual information between the past and future

$$\mathbf{E}(\mathbf{X}) = \lim_{L \rightarrow \infty} I(X_{1:L}; X_{L+1:2L}).$$

It is known that if \mathbf{X} is the output process of a HMM then $\mathbf{E}(\mathbf{X})$ exists [21].

When the alphabet of \mathbf{X} is an ordered alphabet A_n , we define the *permutation excess entropy* of \mathbf{X} [10] by

$$\begin{aligned} \mathbf{E}^*(\mathbf{X}) &= \lim_{L \rightarrow \infty} [H^*(X_{1:L}) - h^*(\mathbf{X})L] \\ &= \sum_{L=1}^{\infty} [H^*(X_L|X_{1:L-1}) - h^*(\mathbf{X})], \end{aligned}$$

if the limit on the right-hand side exists, where $h^*(\mathbf{X}) = \lim_{L \rightarrow \infty} H^*(X_{1:L})/L$ is the *permutation entropy rate* of \mathbf{X} which exists for any finite-alphabet stationary stochastic process and is equal to the entropy rate $h(\mathbf{X})$ [2], [8], [9] and $H^*(X_L|X_{1:L-1}) := H^*(X_{1:L}) - H^*(X_{1:L-1})$.

The following proposition is a generalization of our previous results in [10], [11].

Proposition 5: Let \mathbf{X} be the output process of a HMM $(\Sigma, A_n, \{T^{(a)}\}_{a \in A_n}, \mu)$ with an ergodic internal process. Then, we have

$$\mathbf{E}(\mathbf{X}) = \mathbf{E}^*(\mathbf{X}) = \lim_{L \rightarrow \infty} I^*(X_{1:L}; X_{L+1:2L}),$$

where $I^*(X_{1:L}; X_{L+1:2L}) := H^*(X_{1:L}) + H^*(X_{L+1:2L}) - H^*(X_{1:L}, X_{L+1:2L}) = 2H^*(X_{1:L}) - H^*(X_{1:L}, X_{L+1:2L})$.

Proof: Let $L \geq 1$. We have

$$\begin{aligned} &| [H(X_{1:L}) - h(\mathbf{X})L] - [H^*(X_{1:L}) - h^*(\mathbf{X})L] | \\ &= |H(X_{1:L}) - H^*(X_{1:L})| \\ &\leq \alpha_{\mathbf{X}, L} n \log(L+n) \\ &\leq 2Cn^2 \log(L+n) \gamma^L, \end{aligned}$$

where $C := \max_{x \in A_n} \{C_x\}$, $\gamma := \max_{x \in A_n} \{\gamma_x\} < 1$ and we have used $h(\mathbf{X}) = h^*(\mathbf{X})$ for the first equality, Lemma 4 for the second inequality and Lemma 2 and Lemma 3 for the last inequality. By taking the limit $L \rightarrow \infty$ we obtain $\mathbf{E}(\mathbf{X}) = \mathbf{E}^*(\mathbf{X})$.

To prove

$$\lim_{L \rightarrow \infty} I(X_{1:L}; X_{L+1:2L}) = \lim_{L \rightarrow \infty} I^*(X_{1:L}; X_{L+1:2L}),$$

it is sufficient to show that

$$|H(X_{1:L}, X_{L+1:2L}) - H^*(X_{1:L}, X_{L+1:2L})| \rightarrow 0$$

as $L \rightarrow \infty$. This is because we have

$$\begin{aligned} &|I(X_{1:L}; X_{L+1:2L}) - I^*(X_{1:L}; X_{L+1:2L})| \\ &\leq 2|H(X_{1:L}) - H^*(X_{1:L})| \\ &\quad + |H(X_{1:L}, X_{L+1:2L}) - H^*(X_{1:L}, X_{L+1:2L})|. \end{aligned}$$

However, this can be shown similarly with the above discussion by applying Lemma 4 to the bivariate process $(\mathbf{X}^1, \mathbf{X}^2) := (\mathbf{X}, \mathbf{X})$ and then using Lemma 2 and Lemma 3. ■

C. Transfer Entropy and Momentary Information Transfer

In this subsection we consider two information rates that are measures of coupling direction and strength between two jointly distributed processes and discuss the equalities between them and their permutation versions. One is the rate version

of the transfer entropy [14] and the other is the rate version of the momentary information transfer [16]. Both are particular instances of conditional mutual information [30].

Let (\mathbf{X}, \mathbf{Y}) be a bivariate finite-alphabet stationary stochastic process. We assume that the alphabets of \mathbf{X} and \mathbf{Y} are ordered alphabets A_n and A_m , respectively. For $\tau = 1, 2, \dots$, we define the τ -step transfer entropy rate from \mathbf{Y} to \mathbf{X} by

$$\begin{aligned} t_\tau(\mathbf{Y} \rightarrow \mathbf{X}) &= \lim_{L \rightarrow \infty} [H(X_{1:L+\tau}) - H(X_{1:L}) \\ &\quad - H(X_{1:L+\tau}, Y_{1:L}) + H(X_{1:L}, Y_{1:L})]. \end{aligned}$$

When $\tau = 1$, $t_1(\mathbf{Y} \rightarrow \mathbf{X})$ is called just *transfer entropy rate* [31] from \mathbf{Y} to \mathbf{X} and simply denoted by $t(\mathbf{Y} \rightarrow \mathbf{X})$.

If we introduce the τ -step entropy rate of \mathbf{X} by

$$h_\tau(\mathbf{X}) = \lim_{L \rightarrow \infty} H(X_{L+1:L+\tau} | X_{1:L})$$

and the τ -step conditional entropy rate of \mathbf{X} given \mathbf{Y} by

$$h_\tau(\mathbf{X}|\mathbf{Y}) = \lim_{L \rightarrow \infty} H(X_{L+1:L+\tau} | X_{1:L}, Y_{1:L})$$

then we can write

$$t_\tau(\mathbf{Y} \rightarrow \mathbf{X}) = h_\tau(\mathbf{X}) - h_\tau(\mathbf{X}|\mathbf{Y})$$

because both $h_\tau(\mathbf{X})$ and $h_\tau(\mathbf{X}|\mathbf{Y})$ exist. We call $h_1(\mathbf{X}|\mathbf{Y})$ *conditional entropy rate* and denote it by $h(\mathbf{X}|\mathbf{Y})$ ¹.

$h_\tau(\mathbf{X})$ is *additive*, namely, we always have

$$h_\tau(\mathbf{X}) = \tau h_1(\mathbf{X}) = \tau h(\mathbf{X}).$$

However, for the τ -step conditional entropy rate, the additivity cannot hold in general. It is at most *super-additive*: we only have the inequality

$$h_\tau(\mathbf{X}|\mathbf{Y}) \geq \tau h(\mathbf{X}|\mathbf{Y})$$

in general. Indeed, we have

$$\begin{aligned} h_\tau(\mathbf{X}|\mathbf{Y}) &= \lim_{L \rightarrow \infty} H(X_{L+1:L+\tau} | X_{1:L}, Y_{1:L}) \\ &= \lim_{L \rightarrow \infty} \sum_{\tau'=1}^{\tau} H(X_{L+\tau'} | X_{1:L+\tau'-1}, Y_{1:L}) \\ &\geq \lim_{L \rightarrow \infty} \sum_{\tau'=1}^{\tau} H(X_{L+\tau'} | X_{1:L+\tau'-1}, Y_{1:L+\tau'-1}) \\ &= \tau h(\mathbf{X}|\mathbf{Y}). \end{aligned}$$

This leads to the *sub-additivity* of the τ -step transfer entropy rate:

$$t_\tau(\mathbf{Y} \rightarrow \mathbf{X}) \leq \tau t(\mathbf{Y} \rightarrow \mathbf{X}).$$

An example with the strict inequality can be easily given. Let \mathbf{Y} be an i.i.d. process and \mathbf{X} be defined by $X_1 = Y_1$ and $X_{i+1} = Y_i$. We have $h(\mathbf{X}) = h(\mathbf{Y}) = H(Y_1)$ and $h_\tau(\mathbf{X}|\mathbf{Y}) = (\tau - 1)H(Y_1)$. Hence, $t_\tau(\mathbf{Y} \rightarrow \mathbf{X}) = H(Y_1)$ for any $\tau = 1, 2, \dots$.

There are two permutation versions of the transfer entropy. One is called *symbolic transfer entropy (STE)* [33] and the

¹ Note that the conditional entropy rate here is slightly different from that found in the literature. For example, in [32], conditional entropy rate (called *conditional uncertainty*) is defined by $\lim_{L \rightarrow \infty} H(X_{L+1} | X_{1:L}, Y_{1:L+1})$. The difference from the conditional entropy rate defined here is in whether the conditioning on Y_{L+1} is involved or not.

other is called *transfer entropy on rank vector (TERV)* [34]. Here, we introduce their rate versions as follows: the *rate of STE* from \mathbf{Y} to \mathbf{X} is defined by

$$t_{\tau}^{**}(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} [H^*(X_{1:L}, X_{1+\tau:L+\tau}) - H^*(X_{1:L}) - H^*(X_{1:L}, X_{1+\tau:L+\tau}, Y_{1:L}) + H^*(X_{1:L}, Y_{1:L})]$$

if the limit on the right-hand side exists. The *rate of TERV* from \mathbf{Y} to \mathbf{X} is defined by

$$t_{\tau}^*(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} [H^*(X_{1:L+\tau}) - H^*(X_{1:L}) - H^*(X_{1:L+\tau}, Y_{1:L}) + H^*(X_{1:L}, Y_{1:L})].$$

if the limit on the right-hand side exists. If $\mathbf{E}^*(\mathbf{X})$ exists then, by the definition of the permutation excess entropy, we have

$$h^*(\mathbf{X}) = \lim_{L \rightarrow \infty} [H^*(X_{1:L+1}) - H^*(X_{1:L})].$$

In this case, $t_1^*(\mathbf{Y} \rightarrow \mathbf{X})$ coincides with a quantity called *symbolic transfer entropy rate* introduced in [12].

Proposition 6: Let (\mathbf{X}, \mathbf{Y}) be the output process of a HMM $(\Sigma, A_n \times A_m, \{T^{(a,b)}\}_{(a,b) \in A_n \times A_m}, \mu)$ with an ergodic internal process. Then, we have

$$t_{\tau}(\mathbf{Y} \rightarrow \mathbf{X}) = t_{\tau}^*(\mathbf{Y} \rightarrow \mathbf{X}) = t_{\tau}^{**}(\mathbf{Y} \rightarrow \mathbf{X}).$$

Proof: Since both \mathbf{X} and \mathbf{Y} are the output processes of appropriate HMMs with ergodic internal processes, the equalities follow from the similar discussion with that in the proof of Proposition 5. Indeed, for example, \mathbf{X} is the output process of the HMM $(\Sigma, A_n, \{T^{(a)}\}_{a \in A_n}, \mu)$ where $T^{(a)} := \sum_{b \in A_m} T^{(a,b)}$. ■

A different instance of conditional mutual information called *momentary information transfer* is considered in [16]. It was proposed to improve the ability to detect coupling delays which is lacked in the transfer entropy. Here, we consider its rate version: the *momentary information transfer rate* is defined by

$$m_{\tau}(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} [H(X_{1:L+\tau}, Y_{1:L-1}) - H(X_{1:L+\tau-1}, Y_{1:L-1}) - H(X_{1:L+\tau}, Y_{1:L}) + H(X_{1:L+\tau-1}, Y_{1:L})].$$

Its permutation version called *momentary sorting information transfer rate* is defined by

$$m_{\tau}^*(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} [H^*(X_{1:L+\tau}, Y_{1:L-1}) - H^*(X_{1:L+\tau-1}, Y_{1:L-1}) - H^*(X_{1:L+\tau}, Y_{1:L}) + H^*(X_{1:L+\tau-1}, Y_{1:L})].$$

By the similar discussion with that in the proof of Proposition 6, we obtain the following equality:

Proposition 7: Let (\mathbf{X}, \mathbf{Y}) be the output process of a HMM $(\Sigma, A_n \times A_m, \{T^{(a,b)}\}_{(a,b) \in A_n \times A_m}, \mu)$ with an ergodic internal process. Then, we have

$$m_{\tau}(\mathbf{Y} \rightarrow \mathbf{X}) = m_{\tau}^*(\mathbf{Y} \rightarrow \mathbf{X}).$$

D. Directed Information

Directed information is a measure of coupling direction and strength based on the idea of the *causal conditioning* [18], [35]. Since it is not a particular instance of conditional mutual information, here we treat it separately. In the following presentation, we make use of terminologies from [31], [36].

Let (\mathbf{X}, \mathbf{Y}) be a bivariate finite-alphabet stationary stochastic process. The alphabets of \mathbf{X} and \mathbf{Y} are ordered alphabets A_n and A_m , respectively. The *directed information rate* from \mathbf{Y} to \mathbf{X} is defined by

$$I_{\infty}(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} \frac{1}{L} I(Y_{1:L} \rightarrow X_{1:L})$$

where

$$\begin{aligned} I(Y_{1:L} \rightarrow X_{1:L}) &= \sum_{i=1}^L I(X_i; Y_{1:i} | X_{1:i-1}) \\ &= H(X_{1:L}) - \sum_{i=1}^L H(X_i | X_{1:i-1}, Y_{1:i}). \end{aligned}$$

Note that if $Y_{1:i}$ in the above expression on the right-hand side is replaced by $Y_{1:L}$ then we obtain the mutual information between $X_{1:L}$ and $Y_{1:L}$:

$$I(X_{1:L}; Y_{1:L}) = H(X_{1:L}) - \sum_{i=1}^L H(X_i | X_{1:i-1}, Y_{1:L}).$$

Thus, conditioning on $Y_{1:i}$ for $i = 1, \dots, L$, not on $Y_{1:L}$, distinguishes the directed information from the mutual information. Following [35], we write

$$H(X_{1:L} || Y_{1:L}) := \sum_{i=1}^L H(X_i | X_{1:i-1}, Y_{1:i})$$

and call the quantity *causal conditional entropy*. By using this notation, we have

$$I(Y_{1:L} \rightarrow X_{1:L}) = H(X_{1:L}) - H(X_{1:L} || Y_{1:L}).$$

The permutation version of the directed information rate which we call *symbolic directed information rate* is defined by

$$I_{\infty}^*(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} \frac{1}{L} I^*(Y_{1:L} \rightarrow X_{1:L})$$

if the limit on the right-hand side exists, where

$$\begin{aligned} I^*(Y_{1:L} \rightarrow X_{1:L}) &:= H^*(X_{1:L}) - \sum_{i=1}^L [H^*(X_{1:i}, Y_{1:i}) - H^*(X_{1:i-1}, Y_{1:i})]. \end{aligned}$$

If we write

$$I^*(X_i; Y_{1:i} | X_{1:i-1}) := H^*(X_{1:i}) - H^*(X_{1:i-1}) - H^*(X_{1:i}, Y_{1:i}) + H^*(X_{1:i-1}, Y_{1:i})$$

and

$$H^*(X_{1:L} || Y_{1:L}) := \sum_{i=1}^L [H^*(X_{1:i}, Y_{1:i}) - H^*(X_{1:i-1}, Y_{1:i})]$$

then we have the expressions

$$\begin{aligned} I^*(Y_{1:L} \rightarrow X_{1:L}) &= \sum_{i=1}^L I^*(X_i; Y_{1:i} | X_{1:i-1}) \\ &= H^*(X_{1:L}) - H^*(X_{1:L} || Y_{1:L}). \end{aligned}$$

Proposition 8: Let (\mathbf{X}, \mathbf{Y}) be the output process of a HMM $(\Sigma, A_n \times A_m, \{T^{(a,b)}\}_{(a,b) \in A_n \times A_m}, \mu)$ with an ergodic internal process. Then, we have

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X}) = I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}).$$

Proof: We have

$$\begin{aligned} &|I(Y_{1:L} \rightarrow X_{1:L}) - I^*(Y_{1:L} \rightarrow X_{1:L})| \\ &\leq |H(X_{1:L}) - H^*(X_{1:L})| \\ &\quad + \sum_{i=1}^L |H(X_{1:i}, Y_{1:i}) - H^*(X_{1:i}, Y_{1:i})| \\ &\quad + \sum_{i=1}^L |H(X_{1:i-1}, Y_{1:i}) - H^*(X_{1:i-1}, Y_{1:i})|. \end{aligned}$$

We know that the first term on the right-hand side in the above inequality goes to 0 as $L \rightarrow \infty$. Let us evaluate the second sum. By Lemma 4, it holds that

$$\begin{aligned} &\sum_{i=1}^L |H(X_{1:i}, Y_{1:i}) - H^*(X_{1:i}, Y_{1:i})| \\ &\leq \sum_{i=1}^L (\alpha_{\mathbf{X},i} + \alpha_{\mathbf{Y},i}) [n \log(i+n) + m \log(i+m)] \end{aligned}$$

By Lemma 2 and Lemma 3, we have

$$\sum_{i=1}^L \alpha_{\mathbf{X},i} n \log(i+n) \leq 2Cn^2 \sum_{i=1}^L \gamma^i \log(i+n),$$

where $C := \max_{x \in A_n} \{C_x\}$ and $\gamma := \max_{x \in A_n} \{\gamma_x\} < 1$. It is elementary to show that $\lim_{L \rightarrow \infty} \sum_{i=1}^L \gamma^i \log(i+n)$ is finite. The limits of the other terms are also shown to be finite similarly. Thus, we can conclude that the limit of the second sum is bounded. Similarly, the limit of the third sum is also bounded. The equality in the claim follows immediately. ■

For output processes of HMMs with ergodic internal processes, properties on the directed information rate can be transferred to those on the symbolic directed information rate. Since proofs of them can be given by the same manner as those of the above propositions, here we list some of them without proofs. For the proofs of the properties on the directed information rate, we refer to [35], [36].

Let (\mathbf{X}, \mathbf{Y}) be the output process of a HMM $(\Sigma, A_n \times A_m, \{T^{(a,b)}\}_{(a,b) \in A_n \times A_m}, \mu)$ with an ergodic internal process. Then, we have

(i)

$$I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} I^*(X_L; Y_{1:L} | X_{1:L-1}).$$

This is the permutation version of the equality

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} I(X_L; Y_{1:L} | X_{1:L-1}).$$

(ii)

$$\begin{aligned} I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}) &= I_\infty^*(D\mathbf{Y} \rightarrow \mathbf{X}) \\ &= \lim_{L \rightarrow \infty} I^*(X_L; Y_{1:L-1} | X_{1:L-1}). \end{aligned}$$

Here,

$$I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}) := \lim_{L \rightarrow \infty} \frac{1}{L} I(DY_{1:L} \rightarrow X_{1:L})$$

and

$$I(DY_{1:L} \rightarrow X_{1:L}) := \sum_{i=1}^L I(X_i; Y_{1:i-1} | X_{1:i-1}).$$

The symbol D denotes the one-step delay. $I_\infty^*(D\mathbf{Y} \rightarrow \mathbf{X})$ is the corresponding permutation version. The second equality is the permutation version of the equality

$$I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}) = \lim_{L \rightarrow \infty} I(X_L; Y_{1:L-1} | X_{1:L-1}).$$

Since $I_\infty(D\mathbf{Y} \rightarrow \mathbf{X})$ coincides with the transfer entropy rate, the first equality is just the equality between the transfer entropy rate and the symbolic transfer entropy rate (or the rate of 1-step TERV) proved in Proposition 6 given the second equality.

(iii)

$$\begin{aligned} I_\infty(\mathbf{Y} \rightarrow \mathbf{X} || D\mathbf{Y}) &= I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X} || D\mathbf{Y}) \\ &= \lim_{L \rightarrow \infty} I^*(X_L; Y_L | X_{1:L-1}, Y_{1:L-1}), \end{aligned}$$

where $I_\infty(\mathbf{Y} \rightarrow \mathbf{X} || D\mathbf{Y})$ is called *instantaneous information exchange rate* and is defined by

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X} || D\mathbf{Y}) := \lim_{L \rightarrow \infty} \frac{1}{L} I(Y_{1:L} \rightarrow X_{1:L} || DY_{1:L})$$

and

$$\begin{aligned} &I(Y_{1:L} \rightarrow X_{1:L} || DY_{1:L}) \\ &= H(X_{1:L} || DY_{1:L}) - H(X_{1:L} || Y_{1:L}, DY_{1:L}) \\ &= \sum_{i=1}^L I(X_i; Y_{1:i} | X_{1:i-1}, Y_{1:i-1}) \\ &= \sum_{i=1}^L I(X_i; Y_i | X_{1:i-1}, Y_{1:i-1}). \end{aligned}$$

From the last expression of $I(Y_{1:L} \rightarrow X_{1:L} || DY_{1:L})$, we can obtain

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X} || D\mathbf{Y}) = \lim_{L \rightarrow \infty} I(X_L; Y_L | X_{1:L-1}, Y_{1:L-1}).$$

$I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X} || D\mathbf{Y})$ is the corresponding permutation version and called *symbolic instantaneous information exchange rate*.

(iv)

$$I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}) = I_\infty^*(D\mathbf{Y} \rightarrow \mathbf{X}) + I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X} || D\mathbf{Y}).$$

Namely, the symbolic directed information rate decomposes into the sum of the symbolic transfer entropy rate and the symbolic instantaneous information exchange rate. This follows immediately from (ii), (iii) and the

equality saying that the directed information rate decomposes into the sum of the transfer entropy rate and the instantaneous information exchange rate:

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X}) = I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}) + I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}).$$

(v)

$$I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}) + I_\infty^*(D\mathbf{X} \rightarrow \mathbf{Y}) = I_\infty^*(\mathbf{X}; \mathbf{Y}).$$

This is the permutation version of the equality saying that the mutual information rate between \mathbf{X} and \mathbf{Y} is the sum of the directed information rate from \mathbf{Y} to \mathbf{X} and the transfer entropy rate from \mathbf{X} to \mathbf{Y} :

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X}) + I_\infty(D\mathbf{X} \rightarrow \mathbf{Y}) = I_\infty(\mathbf{X}; \mathbf{Y}),$$

where

$$I_\infty(\mathbf{X}; \mathbf{Y}) := \lim_{L \rightarrow \infty} \frac{1}{L} I(X_{1:L}; Y_{1:L})$$

is the *mutual information rate* and $I_\infty^*(\mathbf{X}; \mathbf{Y})$ is its permutation version called *symbolic mutual information rate*. It is known that they are equal for any bivariate finite-alphabet stationary stochastic process [12]. Thus, the symbolic mutual information rate between \mathbf{X} and \mathbf{Y} is the sum of the symbolic directed information rate from \mathbf{Y} to \mathbf{X} and the symbolic transfer entropy rate from \mathbf{X} to \mathbf{Y} .

We can also introduce the permutation version of the *causal conditional directed information rate* and prove the corresponding properties. To be precise, let us consider a multivariate finite-alphabet stationary stochastic process $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}^1, \dots, \mathbf{Z}^k)$ with the alphabet $A_n \times A_m \times A_{l_1} \times \dots \times A_{l_k}$. The *causal conditional directed information rate* from \mathbf{Y} to \mathbf{X} given $\mathcal{Z} := (\mathbf{Z}^1, \dots, \mathbf{Z}^k)$ is defined by

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) := \lim_{L \rightarrow \infty} \frac{1}{L} I(Y_{1:L} \rightarrow X_{1:L}||Z_{1:L}^1, \dots, Z_{1:L}^k) \quad (\text{iv}')$$

where

$$\begin{aligned} & I(Y_{1:L} \rightarrow X_{1:L}||Z_{1:L}^1, \dots, Z_{1:L}^k) \\ &= H(X_{1:L}||Z_{1:L}^1, \dots, Z_{1:L}^k) - H(X_{1:L}||Y_{1:L}, Z_{1:L}^1, \dots, Z_{1:L}^k) \\ &= \sum_{i=1}^L I(X_i; Y_{1:i}|X_{1:i-1}, Z_{1:L}^1, \dots, Z_{1:L}^k). \end{aligned}$$

Corresponding to Proposition 8, we have the following equality if $(\mathbf{X}, \mathbf{Y}, \mathcal{Z})$ is the output process of a HMM with an ergodic internal process:

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) = I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}),$$

where $I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z})$ is the *symbolic causal conditional directed information rate* which is defined by the same manner as the symbolic directed information rate. The following properties also hold: assume that $(\mathbf{X}, \mathbf{Y}, \mathcal{Z})$ is the output process of a HMM with an ergodic internal process. Then, we have

(i')

$$\begin{aligned} & I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) \\ &= \lim_{L \rightarrow \infty} I^*(X_L; Y_{1:L}|X_{1:L-1}, Z_{1:L}^1, \dots, Z_{1:L}^k). \end{aligned}$$

This is the permutation version of the equality

$$\begin{aligned} & I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) \\ &= \lim_{L \rightarrow \infty} I(X_L; Y_{1:L}|X_{1:L-1}, Z_{1:L}^1, \dots, Z_{1:L}^k). \end{aligned}$$

(ii')

$$\begin{aligned} & I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) = I_\infty^*(D\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) \\ &= \lim_{L \rightarrow \infty} I^*(X_L; Y_{1:L-1}|X_{1:L-1}, Z_{1:L}^1, \dots, Z_{1:L}^k). \end{aligned}$$

The second equality is the permutation version of the equality

$$\begin{aligned} & I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) \\ &= \lim_{L \rightarrow \infty} I(X_L; Y_{1:L-1}|X_{1:L-1}, Z_{1:L}^1, \dots, Z_{1:L}^k). \end{aligned}$$

The quantities $I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z})$ and $I_\infty^*(D\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z})$ are called *causal conditional transfer entropy rate* and *symbolic causal conditional transfer entropy rate*, respectively.

$$\begin{aligned} & I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}, \mathcal{Z}) = I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}, \mathcal{Z}) \\ &= \lim_{L \rightarrow \infty} I^*(X_L; Y_L|X_{1:L-1}, Y_{1:L-1}, Z_{1:L}^1, \dots, Z_{1:L}^k), \end{aligned}$$

where $I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}, \mathcal{Z})$ is called *causal conditional instantaneous information exchange rate*. The second equality is the permutation version of the equality

$$\begin{aligned} & I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}, \mathcal{Z}) \\ &= \lim_{L \rightarrow \infty} I(X_L; Y_L|X_{1:L-1}, Y_{1:L-1}, Z_{1:L}^1, \dots, Z_{1:L}^k). \end{aligned}$$

$I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}, \mathcal{Z})$ is the permutation version and is called *symbolic causal conditional instantaneous information exchange rate*.

$$\begin{aligned} & I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) \\ &= I_\infty^*(D\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) + I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}, \mathcal{Z}). \end{aligned}$$

This is the permutation version of the equality

$$\begin{aligned} & I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) \\ &= I_\infty(D\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) + I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||D\mathbf{Y}, \mathcal{Z}). \end{aligned}$$

(v')

$$I_\infty^*(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) + I_\infty^*(D\mathbf{X} \rightarrow \mathbf{Y}||\mathcal{Z}) = I_\infty^*(\mathbf{X}; \mathbf{Y}||\mathcal{Z}).$$

This is the permutation version of the equality

$$I_\infty(\mathbf{Y} \rightarrow \mathbf{X}||\mathcal{Z}) + I_\infty(D\mathbf{X} \rightarrow \mathbf{Y}||\mathcal{Z}) = I_\infty(\mathbf{X}; \mathbf{Y}||\mathcal{Z}),$$

where

$$\begin{aligned} & I_\infty(\mathbf{X}; \mathbf{Y}||\mathcal{Z}) \\ &:= \lim_{L \rightarrow \infty} \frac{1}{L} [H(X_{1:L}||Z_{1:L}^1, \dots, Z_{1:L}^k) \\ &\quad + H(Y_{1:L}||Z_{1:L}^1, \dots, Z_{1:L}^k) \\ &\quad - H(X_{1:L}, Y_{1:L}||Z_{1:L}^1, \dots, Z_{1:L}^k)] \end{aligned}$$

is the *causal conditional mutual information rate* and $I_\infty^*(\mathbf{X}; \mathbf{Y}||\mathcal{Z})$ is its permutation version called *symbolic*

causal conditional mutual information rate. It can be shown that

$$I_\infty(\mathbf{X}; \mathbf{Y} | \mathcal{Z}) = I_\infty^*(\mathbf{X}; \mathbf{Y} | \mathcal{Z})$$

if $(\mathbf{X}, \mathbf{Y}, \mathcal{Z})$ is the output process of a HMM with an ergodic internal process.

V. DISCUSSION

In this section, we discuss how our theoretical results in this paper are related to the previous work in the literature.

Being confronted with real time series data, we cannot take the limit of large length of words. Hence, we have to estimate information rates with finite length of words. In such situation, one permutation method could have some advantages to the other permutation methods. As a matter of fact, TERV was originally proposed as an improved version of STE [34]. However, it has been unclear whether they coincide in the limit of large length of permutations. In this paper, we provide a partial answer to this question: the two permutation versions of the transfer entropy rate, the rate of STE and the rate of TERV, are equivalent to the transfer entropy rate for bivariate processes generated by HMMs with ergodic internal processes.

Granger causality graph [37] is a model of causal dependence structure in multivariate stationary stochastic processes. Given a multivariate stationary stochastic process, nodes in a Granger causality graph are components of the process. There are two types of edges: one is directed and the other is undirected. The absence of a directed edge from one node to another node indicates the lack of the Granger cause from the former to the latter relative to the other remaining processes. Similarly, the absence of an undirected edge between two nodes indicates the lack of the instantaneous cause between them relative to the other remaining processes. Amblard and Michel [31], [36] proposed that the Granger causality graph can be constructed based on the directed information theory: let $\mathcal{X} = (\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^m)$ be a multivariate finite-alphabet stationary stochastic process with the alphabet $A_{n_1} \times A_{n_2} \times \dots \times A_{n_m}$ and (V, E_d, E_u) be the Granger causality graph of the process \mathcal{X} where $V = \{1, 2, \dots, m\}$ is the set of nodes, E_d is the set of directed edges and E_u is the set of undirected edges. Their proposal is that

- (i) for any $i, j \in V$, $(i, j) \notin E_d$ if and only if $I_\infty(D\mathbf{X}^i \rightarrow \mathbf{X}^j | \mathcal{X} \setminus \{\mathbf{X}^i, \mathbf{X}^j\}) = 0$,
- (ii) for any $i, j \in V$, $(i, j) \notin E_u$ if and only if $I_\infty(\mathbf{X}^i \rightarrow \mathbf{X}^j | D\mathbf{X}^i, \mathcal{X} \setminus \{\mathbf{X}^i, \mathbf{X}^j\}) = 0$.

Thus, in the Granger causality graph construction proposed in [31], [36], the causal conditional transfer entropy rate captures the Granger cause from one process to another process relative to the other remaining processes. On the other hand, the causal conditional instantaneous information exchange rate captures the instantaneous cause between two processes relative to the other remaining processes.

Now, let us consider the case when \mathcal{X} is the output process of a HMM with an ergodic internal process. Then, from the results of Section IV-D, we have

- (i') for any $i, j \in V$, $(i, j) \notin E_d$ if and only if $I_\infty^*(D\mathbf{X}^i \rightarrow \mathbf{X}^j | \mathcal{X} \setminus \{\mathbf{X}^i, \mathbf{X}^j\}) = 0$,

- (ii') for any $i, j \in V$, $(i, j) \notin E_u$ if and only if $I_\infty^*(\mathbf{X}^i \rightarrow \mathbf{X}^j | D\mathbf{X}^i, \mathcal{X} \setminus \{\mathbf{X}^i, \mathbf{X}^j\}) = 0$.

Thus, the Granger causality graphs in the sense of [31], [36] for multivariate processes generated by HMMs with ergodic internal processes can be captured by the language of the permutation entropy: the symbolic causal conditional transfer entropy rate and the symbolic instantaneous information exchange rate. This statement opens up a possibility of the permutation approach to the problem of assessing the causal dependence structure of multivariate stationary stochastic processes. However, of course, the details of the practical implementation should be an issue of further study.

ACKNOWLEDGMENT

The authors would like to thank D. Kugiumtzis for his useful comments and discussion on the relationship between STE and TERV. TH was supported by the JST PRESTO program.

REFERENCES

- [1] C. Bandt and B. Pompe, "Permutation entropy: a natural complexity measure for time series," *Phys. Rev. Lett.*, vol. 88, p. 174102, 2002.
- [2] J. M. Amigó, *Permutation Complexity in Dynamical Systems*. Springer-Verlag Berlin Heidelberg, 2010.
- [3] A. Bahraminasab, F. Ghasemi, A. Stefanovska, P. V. E. McClintock, and H. Kantz, "Direction of coupling from phases of interacting oscillators: A permutation information approach," *Phys. Rev. Lett.*, vol. 100, p. 084101, 2008.
- [4] Y. H. Cao, W. W. Tung, J. B. Gao, V. A. Protopopescu, and L. M. Hively, "Detecting dynamical changes in time series using the permutation entropy," *Phys. Rev. E*, vol. 70, p. 046217, 2004.
- [5] O. A. Rosso, H. A. Larrondo, M. T. Martin, A. Plastino, and M. A. Fuentes, "Distinguishing noise from chaos," *Phys. Rev. Lett.*, vol. 99, p. 154102, 2007.
- [6] C. Bandt, G. Keller, and B. Pompe, "Entropy of interval maps via permutations," *Nonlinearity*, vol. 15, pp. 1595–1602, 2002.
- [7] K. Keller and M. Sinn, "Kolmogorov-sinai entropy from the ordinal viewpoint," *Physica D*, vol. 239, pp. 997–1000, 2010.
- [8] J. M. Amigó, M. B. Kennel, and L. Kocarev, "The permutation entropy rate equals the metric entropy rate for ergodic information sources and ergodic dynamical systems," *Physica D*, vol. 210, pp. 77–95, 2005.
- [9] J. M. Amigó, "The equality of kolmogorov-sinai entropy and metric permutation entropy generalized," *Physica D*, vol. 241, pp. 789–793, 2012.
- [10] T. Haruna and K. Nakajima, "Permutation complexity via duality between values and orderings," *Physica D*, vol. 240, pp. 1370–1377, 2011.
- [11] —, "Permutation excess entropy and mutual information between the past and future," arXiv:1112.2491v1.
- [12] —, "Symbolic transfer entropy rate is equal to transfer entropy rate for bivariate finite-alphabet stationary ergodic markov processes," arXiv:1112.2493v2.
- [13] J. P. Crutchfield and D. P. Feldman, "Regularities unseen, randomness observed: Levels of entropy convergence," *Chaos*, vol. 15, pp. 25–54, 2003.
- [14] T. Schreiber, "Measuring information transfer," *Phys. Rev. Lett.*, vol. 85, pp. 461–464, 2000.
- [15] A. Kaiser and T. Schreiber, "Information transfer in continuous processes," *Physica D*, vol. 166, pp. 43–62, 2002.
- [16] B. Pompe and J. Runge, "Momentary information transfer as a coupling measure of time series," *Phys. Rev. E*, vol. 83, p. 051122, 2011.
- [17] H. Marko, "The bidirectional communication theory—a generalization of information theory," *IEEE Transactions on Communications*, vol. 21, pp. 1345–1351, 1973.
- [18] J. L. Massey, "Causality, feedback and directed information," in *Proc. Intl. Symp. on Infor. Th. and its Applications*. Waikiki, Hawaii, 1990.
- [19] D. R. Upper, *Theory and Algorithms for Hidden Markov Models and Generalized Hidden Markov Models*. PhD thesis, University of California, Berkeley, 1997.
- [20] P. Walters, *An Introduction to Ergodic Theory*. Springer-Verlag, New York, 1982.

- [21] W. Löhner, *Models of Discrete Time Stochastic Processes and Associated Complexity Measures*. PhD thesis, Max Planck Institute for Mathematics in the Sciences, Leipzig, 2010.
- [22] A. Katok and B. Hasselblatt, *Introduction to the Modern Theory of Dynamical Systems*. Cambridge University Press, 1995.
- [23] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. John Wiley & Sons, Inc, 1991.
- [24] D. V. Arnold, "Information-theoretic analysis of phase transitions," *Complex Systems*, vol. 10, pp. 143–155, 1996.
- [25] W. Bialek, I. Nemenman, and N. Tishby, "Predictability, complexity, and learning," *Neural Computation*, vol. 13, pp. 2409–2463, 2001.
- [26] D. P. Feldman, C. S. McTague, and J. P. Crutchfield, "The organization of intrinsic computation: complexity-entropy diagrams and the diversity of natural information processing," *Chaos*, vol. 18, p. 043106, 2008.
- [27] P. Grassberger, "Toward a quantitative theory of self-generated complexity," *Int. J. Theor. Phys.*, vol. 25, pp. 907–938, 1986.
- [28] W. Li, "On the relationship between complexity and entropy for markov chains and regular languages," *Complex Systems*, vol. 5, pp. 381–399, 1991.
- [29] R. Shaw, *The Dripping Faucet as a Model Chaotic System*. Aerial Press, Santa Cruz, California, 1984.
- [30] S. Frenzel and B. Pompe, "Partial mutual information for coupling analysis of multivariate time series," *Phys. Rev. Lett.*, vol. 99, p. 204101, 2007.
- [31] P.-O. Amblard and O. J. J. Michel, "On directed information theory and granger causality graphs," *J. Comput. Neurosci.*, vol. 30, pp. 7–16, 2011.
- [32] R. Ash, *Information Theory*. Wiley Interscience, New York, 1965.
- [33] M. Staniek and K. Lehnertz, "Symbolic transfer entropy," *Phys. Rev. Lett.*, vol. 100, p. 158101, 2008.
- [34] D. Kugiumtzis, "Transfer entropy on rank vectors," arXiv:1007.0357v1.
- [35] G. Kramer, *Directed information for channels with feedback*. PhD thesis, Swiss Federal Institute of Technology, Zurich, 1998.
- [36] P.-O. Amblard and O. J. J. Michel, "Relating granger causality to directed information theory for networks of stochastic processes," arXiv:0911.2873v4.
- [37] R. Dahlaus and M. Eichler, "Causality and graphical models in time series analysis," in *Highly structured stochastic systems*, P. Green, N. Hjort, and S. Richardson, Eds. Oxford University Press, 2003, pp. 115–137.